

Aidan Ewart

aidanprattewart 'at' gmail 'dot' com

github.com/Baidicoot

linkedin.com/in/aidanewart

SUMMARY

I am an undergraduate student studying Mathematics at the University of Bristol. I am an experienced programmer (8+ years) and am proficient at a number of programming languages. Currently, my interests are primarily in AI safety, and I have authored a number of papers in that field.

SKILLS

Programming Languages

Python, Haskell, Rust, C, C++, JavaScript, Coq

Frameworks

PyTorch, HuggingFace Accelerate

Other Software

Git, Windows, Unix, TeX, SSH, WebDev

EDUCATION

University of Bristol

Sept 2022 - July 2026

Undergraduate Mathematics (current 3rd year)

1st in all graded classes

Royal Grammar School

Sept 2019 - July 2022

A-Levels in Maths, Further Maths, Physics, Computer Science

*A*A*A*A**

EXPERIENCE

Haize Labs

July 2024 - Oct 2024

Research Intern

New York, NY

- Early stage startup, interned when they had <5 employees
- Built up their automated red-teaming infrastructure, developing novel automated attacks
- Contributed to Haize's contract work with Anthropic and red-teaming of new OpenAI models

MATS 5.0

Jan 2024 - Mar 2024

Scholar

Berkeley, CA

- Developed a training method (T-LAT) for unlearning and adversarial robustness in language models
- T-LAT was a Pareto improvement over prior robustness methods against automated attacks
- Demonstrated applications in removing backdoors inserted via data poisoning

SELECT PUBLICATIONS

* = *order randomized/equal contribution*

Sparse Autoencoders Find Highly Interpretable Features in Language Models

ICLR Conference Paper

H. Cunningham, A. Ewart*, L. Riggs*, R. Huben, L. Sharkey*

- Conference paper at ICLR 2024
- Workshop paper at ATTRIB @ NeurIPS 202

Targeted Latent Adversarial Training Improves Robustness to Persistent Harmful Behaviors in LLMs

Preprint

A. Ewart, A. Sheshadri*, P. Guo, A. Lynch, C. Wu, V. Hebbar, H. Sleight, A. Cooper Stickland, E. Perez, D. Hadfield-Menell, S. Casper*

Eight Methods to Evaluate Robust Unlearning in LLMs

Preprint

A. Lynch, P. Guo*, A. Ewart*, S. Casper, D. Hadfield-Menell*

SELECT PROJECTS

Functional Programming Language Compiler

Source on GitHub

Haskell, x86 Assembly, C

- Implemented a compiler for a Lisp-like high-level programming language
- Frontend includes Hindley-Milner typechecking and inference, a module/imports system, compilation with continuations
- Backend includes program optimisation, register allocation, compilation to x86 assembly and C

Proof Assistant

Source on GitHub

Lua, Haskell

- Implemented a theorem-proving DSL for Lua
- Proof assistant uses a Martin-Löf style type system complete with type inference via unification
- Includes a customisable notation system in the style of Coq